

CH5350: Applied Time-Series Analysis

Fisher's information and Properties of Estimators

Arun K. Tangirala

Department of Chemical Engineering, IIT Madras

Learning Goals

In this lecture, we shall learn the following concepts / topics:

- ▶ Goodness of estimators
- ▶ Fisher information
- ▶ Bias and Variance
- ▶ Efficiency and C-R Inequality
- ▶ Mean Square Error and MMSE
- ▶ Consistency
- ▶ Distribution of estimates

Goodness of an estimate

The quality of an estimate critically depends on two factors:

- 1 **Quality of data**, i.e., how “informative” it is w.r.t. unknowns. A quantitative measure of information is required for this purpose.
- 2 **Goodness of estimator**, i.e., how “well” the estimator is able to process the information to extract the unknown. Of particular interest are the efficiency (related to variance) and consistency (concerned with convergence of estimates) properties.

Metrics characterize the estimator; however, remember that **a fundamental requirement for obtaining a good estimate is that the data should be informative** (with respect to the parameters).

Fisher information

Fisher introduced the notion of information in a data through a series of works by and some existing results. Intuitively, larger the information index is, the “better” the estimator is.

The Fisher information (FI) (Fisher, 1922, 1950) is based on the **likelihood function** of the given data.

The likelihood function stems from the notion of conditional probability, i.e., the probability of observing an event within the vicinity of given data.

Likelihood function

The probability of obtaining data within the vicinity of \mathbf{y}_N is given by (with some abuse of notation)

$$\Pr(\mathbf{y}_N < \mathbf{Y} < \mathbf{y}_N + d\mathbf{y}_N) = f(\mathbf{y}_N|\boldsymbol{\theta})d\mathbf{y}_N \propto f(\mathbf{y}_N|\boldsymbol{\theta}) \quad (1)$$

For a given \mathbf{y}_N , the probability is solely a function of $\boldsymbol{\theta}$. Fisher's argument (and the likes of it) rests on the **maximum likelihood** premise that

Among all possible values of $\boldsymbol{\theta}$, the one that maximizes the probability, i.e., the one that renders the event most likely is the winner!

Likelihood function

The likelihood function (of θ) is, therefore (for continuous RVs), defined as

$$\boxed{l(\theta, \mathbf{y}) = f(\mathbf{y}; \theta)} \quad (\text{or } f(\mathbf{y}|\theta)) \quad (2)$$

where \mathbf{y} is the vector of N observations.

- ▶ The fundamental difference between $l(\theta|\mathbf{y})$ and $f(\mathbf{y}|\theta)$ is that the former is a function of a *deterministic* vector θ , while the latter is a function of the *random* vector \mathbf{y} (given θ).
- ▶ Likelihood function belongs to the world of **statistics** while the p.d.f. belongs to the world of **probability**!

Fisher's information quantifies “how informative” a vector of observations is about a parameter θ (or $\boldsymbol{\theta}$). It rests on the following quantities (assume **single parameter**):

$$l(\theta, \mathbf{y}) = f(\mathbf{y}; \theta) \text{ (or } f(\mathbf{y}|\theta)) \quad \text{(likelihood function)} \quad (3)$$

$$L(\theta, \mathbf{y}) = \ln l(\theta, \mathbf{y}) \quad \text{(log-likelihood function)} \quad (4)$$

$$S(\theta; \mathbf{y}) = \frac{\partial}{\partial \theta} \ln f(\mathbf{y}; \theta) = \frac{\partial}{\partial \theta} L(\theta, \mathbf{y}) \quad \text{(score function)} \quad (5)$$

where \mathbf{y} is the set of observations and θ is the parameter to be estimated.

Further **assume that the p.d.f. is regular** \implies (i) $\partial L / \partial \theta$ exists and is finite and (ii) the operations of integration w.r.t. y and differentiation w.r.t. θ can be interchanged.

FI measures the variability in sensitivity of likelihood, i.e., the score function, across the outcome space (of \mathbf{y}).

The **Fisher information** of a parameter θ in \mathbf{y} is defined as

$$I(\theta) = \text{var}(S) = E \left(\left(\frac{\partial L}{\partial \theta} \right)^2 \right) \quad (6)$$

Under the regularity assumption, it can be shown that

$$\mu_S = E(S|\theta) = 0, \quad \text{var}(S|\theta) = E(S^2) = E \left(\left(\frac{\partial L(\mathbf{y}, \theta)}{\partial \theta} \right)^2 \right) \quad (7)$$

Since

$$E \left(\left(\frac{\partial L}{\partial \theta} \right)^2 \right) = -E \left(\frac{\partial^2 L}{\partial \theta^2} \right) \quad (8)$$

the information can also be computed as

$$I(\theta) = -E \left(\frac{\partial^2 L}{\partial \theta^2} \right) = -E \left(\frac{\partial S}{\partial \theta} \right) \quad (9)$$

Example 1: Information about mean and variance

Consider the case of estimating mean μ and variance σ^2 of a random signal.

Mean and variance

Given that a stationary signal $y[k] \sim \mathcal{N}(\mu, \sigma^2)$, determine (i) $I(\mu)$ and (ii) $I(\sigma^2)$ in a single observation.

- 1 The log-likelihood function (assuming σ^2 is known) is

$$L(\mu; Y) = \ln f(y|\mu) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2} \frac{(y - \mu)^2}{\sigma^2} \quad (10)$$

Example 1

... contd.

The Fisher information on $\theta = \mu$ using (9) is then

$$I(\mu) = -E \left(\frac{\partial^2 L}{\partial \theta^2} \right) = \frac{1}{\sigma^2} \quad (11)$$

Thus, we have a meaningful result. As the variance (spread of possible outcomes) decreases, the information on μ in a *single sample* increases.

Example 1

... contd.

② Now, $\theta = \sigma^2$. The information contained in a single observation is

$$I(\sigma^2) = -E \left(\frac{\partial^2 L}{\partial \theta^2} \right) = -E \left(\frac{1}{2\sigma^4} - \frac{(y - \mu)^2}{\sigma^6} \right) = \frac{1}{2\sigma^4} \quad (12)$$

Example 1

... contd.

- ③ On the other hand, if the parameter of interest is the standard deviation $\theta = \sigma$, the information contained is

$$I(\sigma) = -E \left(\frac{\partial^2 L}{\partial \theta^2} \right) = -E \left(\frac{1}{\sigma^2} - 3 \frac{(y - \mu)^2}{\sigma^4} \right) = \frac{2}{\sigma^2} \quad (13)$$

Thus, $I(\sigma^2) \neq (I(\sigma))^2$. The *information is not commutative with respect to a functional of the parameter $\phi(\theta)$* .

In general, the FI $I(\phi(\theta))$ is related to $I(\theta)$ through

$$\boxed{I(\theta) = \left(\frac{d\phi}{d\theta} \right)^2 I(\phi(\theta))} \quad (14)$$

Fisher information: General case

Generalizing (9) to the case of $p \times 1$ parameter vector $\boldsymbol{\theta}$ contained in N observations, the **information matrix** results:

$$\mathbf{I}_{ij}(\boldsymbol{\theta}) = \text{cov}(S_i, S_j) = E(S_i(\mathbf{Y}_N)S_j(\mathbf{Y}_N)) = -E\left(\frac{\partial^2}{\partial\theta_i\partial\theta_j}L(\boldsymbol{\theta}; \mathbf{y}_N)\right) \quad i, j = 1, \dots, p$$
(15)

where S_i is the i^{th} score statistic,

$$S_i = \frac{\partial}{\partial\theta_i} \ln f(Y_N|\boldsymbol{\theta})$$
(16)

where $f(Y_n|\boldsymbol{\theta})$ is the joint p.d.f. of the N observations \mathbf{y} .

Think: What do the off-diagonal elements signify?

Example 2: Estimating μ , σ^2 from N observations

Information in N observations

Compute the information contained in N samples of a GWN process $y[k] \sim \mathcal{N}(\mu, \sigma^2)$ w.r.t.: (i) $\theta = \mu$ and σ^2 is known, (ii) $\theta = \sigma^2$ and (iii) $\theta = \begin{bmatrix} \mu & \sigma^2 \end{bmatrix}^T$.

Solution: For all the three cases,

$$f(\mathbf{Y}_N | (\mu, \sigma^2)) = \prod_{k=0}^{N-1} \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y[k] - \mu)^2}{2\sigma^2}\right)$$

Information in N observations

- 1 Constructing the log-likelihood from $f(\theta; \mathbf{y}_N)$ gives

$$S(\theta; \mathbf{y}_N) = \frac{\sum_{k=0}^{N-1} (y[k] - \mu)}{\sigma^2}$$

$$\text{Applying (15), } I(\mu) = -E \left(\frac{\partial S}{\partial \theta} \right) = \frac{N}{\sigma^2}$$

Example 2

... contd.

② For this case, $S(\theta; \mathbf{y}_N) = -\frac{N}{2\sigma^2} + \frac{1}{2\sigma^2} \sum_{k=0}^{N-1} (y[k] - \mu)^2$

Applying (15), $I(\sigma^2) = -\frac{\partial S}{\partial \theta} = \frac{N}{2\sigma^4}$

③ Denote $\theta_1 = \mu$ and $\theta_2 = \sigma^2$, $\boldsymbol{\theta} = \begin{bmatrix} \theta_1 & \theta_2 \end{bmatrix}^T$ The log-likelihood function is

$$L(\boldsymbol{\theta}; \mathbf{y}_N) = c - \frac{N}{2} \ln \theta_2 - \frac{1}{2\theta_2} \sum_{k=0}^{N-1} (y[k] - \theta_1)^2 \quad (17)$$

Example 2

... contd.

The information matrix is thus

$$\mathbf{I}(\boldsymbol{\theta}) = -E \left(\begin{bmatrix} \frac{\partial^2 L}{\partial \theta_1^2} & \frac{\partial^2 L}{\partial \theta_1 \partial \theta_2} \\ \frac{\partial^2 L}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 L}{\partial \theta_2^2} \end{bmatrix} \right) = \begin{bmatrix} \frac{N}{\sigma^2} & 0 \\ 0 & \frac{N}{2\sigma^4} \end{bmatrix} \quad (18)$$

Thus, the estimates of mean and variance of a WN process do not affect each other, *i.e.*, these parameters can be estimated individually.

Remarks

- ▶ The Fisher information is a localized version (in the parameter space) of the more general **Kullback-Leibler information** (KLI) in the vicinity of the true parameters. The KLI measures the information loss incurred in approximating a true probability distribution with a model distribution.
- ▶ *Information* is leveraged on two factors: (i) the number and type of unknown(s) that have to be estimated and (ii) how these unknown(s) enter the *model*. Implications of these results are felt in model estimation and in input design.
- ▶ From the examples, we learn that by increasing the sample size, the increase in information is proportional. However, this is not the case when the observations are correlated. In fact, for that case $I_N(\theta) < NI_1(\theta)$.

Bias

One of the foremost expectations of an estimator is that it gives **accurate** estimates.

Definition

An estimator $\hat{\theta}$ is said to be *accurate* or *unbiased* if and only if

$$\mu_{\hat{\theta}} = E(\hat{\theta}) = \theta_0 \quad (19)$$

In plain language, the average of estimates across the records should yield the true value.

The difference $\triangle\hat{\theta} = E(\hat{\theta}) - \theta_0$ is said to be the **bias** of that estimator.

Example: Unbiased estimator

Example

The sample mean estimator $\bar{y} = \frac{1}{N} \sum_{k=0}^{N-1} y[k]$ is unbiased since

$$E(\bar{y}) = E\left(\frac{1}{N} \sum_{k=0}^{N-1} y[k]\right) = \frac{1}{N} \sum_{k=0}^{N-1} E(y[k]) = \mu \quad (20)$$

assuming $y[k]$ to be **stationary**.

Remarks

- ▶ The averaging in (19) is across all possible records of data and *not* along time. From this viewpoint, the definition has limited practical value since it is extremely rare to obtain multiple records of data.
- ▶ Unbiasedness is nevertheless a useful requirement for comparing performance of two estimators.

An unbiased estimator is desirable. However, what is generally more important is the **spread** of estimates when different realizations are presented. This is measured by the **variance** of the estimator.

Variance of estimators

Definition

The variance of an estimator (estimate) is defined as

$$\sigma_{\hat{\theta}}^2 = E((\hat{\theta} - \mu_{\hat{\theta}})^2) = E((\hat{\theta} - E(\hat{\theta}))^2) \quad (21)$$

- ▶ Observe that the definition is w.r.t the average of the estimator, $\mu_{\hat{\theta}}$ and *not* with respect to its true value, θ_0 . When the estimator is unbiased, $E(\hat{\theta}) = \theta_0$.
- ▶ The square root of the variance in (21) is the **standard error** in an estimate.
- ▶ It is obviously desirable to have the variance of estimate much lower than that in the data itself.

Remarks

The variance expression is useful in many different ways:

- i. Computation of error in estimates.
- ii. Constructing confidence regions for the true parameters.
- iii. Design of experiments, i.e., knowing how experimental factors can be adjusted to achieve more reliable (precise) estimates.

Example: Variance

Variance of sample mean

Using Definition 3,

$$\begin{aligned}\sigma_{\bar{y}}^2 &= E((\bar{y} - E(\bar{y}))^2) = E\left(\left(\frac{1}{N} \sum_{k=0}^{N-1} y[k] - \mu_y\right)^2\right) = E\left(\left(\frac{1}{N} \sum_{k=1}^N (y[k] - \mu_y)\right)^2\right) \\&= \frac{1}{N^2} E\left(\sum_{k=1}^N (y[k] - \mu_y)^2\right) + \frac{1}{N^2} E\left(\sum_{n=1}^N \sum_{m=1, m \neq n}^N (y[n] - \mu_y)(y[m] - \mu_y)\right) \\&= \frac{1}{N^2} \left(\sum_{k=1}^N E(y[k] - \mu_y)^2\right) + \frac{1}{N^2} \left(\sum_{n=1}^N \sum_{m=1, m \neq n}^N E(y[n] - \mu_{y,n})(y[m] - \mu_{y,n})\right)\end{aligned}$$

Example

... contd.

The summand in the second term can be easily recognized as the ACVF $y[k]$.

When the signal is WN, *i.e.*,

$$y[k] = c + e[k], \quad e[k] \sim \text{GWN}(0, \sigma_e^2)$$

the variability of sample mean is

$$\boxed{\sigma_y^2 = \frac{\sigma_y^2}{N} = \frac{\sigma_e^2}{N}} \quad (22)$$

Remarks

- ▶ $\text{var}(\bar{y}) \propto \sigma_y^2$ (for a fixed sample size). Intuitively this is a meaningful result. However, we have no control over this factor.
- ▶ $\text{var}(\bar{y}) \propto 1/N, \implies \sigma_{\bar{y}}^2 \rightarrow 0$ as $N \rightarrow \infty$. This is an interesting result and also a good feature of the estimator. Thus, we are able to shrink the variability (and the error) in the estimate by collecting more samples.
- ▶ As we shall shortly learn, (unbiased) estimators that possess this feature are known to be **consistent**.
- ▶ The true mean has no bearing on the variability (of the sample mean), which is again a sensible result.

Variance of vector of parameters

When $\hat{\boldsymbol{\theta}}$ is a $p \times 1$ vector, we have a variance-covariance **matrix**,

$$\begin{aligned} \text{Var}(\hat{\boldsymbol{\theta}}) = \Sigma_{\hat{\boldsymbol{\theta}}} &= E((\hat{\boldsymbol{\theta}} - E(\hat{\boldsymbol{\theta}}))(\hat{\boldsymbol{\theta}} - E(\hat{\boldsymbol{\theta}}))^T) \\ &= \begin{bmatrix} \sigma_{\hat{\theta}_1}^2 & \sigma_{\hat{\theta}_1\hat{\theta}_2} & \cdots & \sigma_{\hat{\theta}_1\hat{\theta}_p} \\ \sigma_{\hat{\theta}_2\hat{\theta}_1} & \sigma_{\hat{\theta}_2}^2 & \cdots & \sigma_{\hat{\theta}_2\hat{\theta}_p} \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{\hat{\theta}_p\hat{\theta}_1} & \sigma_{\hat{\theta}_p\hat{\theta}_2} & \vdots & \sigma_{\hat{\theta}_p}^2 \end{bmatrix} \end{aligned} \quad (23)$$

Remarks

- ▶ It is a symmetric matrix with the off-diagonal elements reflecting the error incurred in estimating a pair of parameters jointly.
- ▶ **A diagonal $\Sigma_{\hat{\theta}}$ connotes that the parameters can be estimated on an individual basis. In practice, the $\text{trace}(\Sigma_{\hat{\theta}})$ and the diagonal elements of $\Sigma_{\hat{\theta}}$ find wider utility.**

Minimum Variance Unbiased Estimator (MVUE)

Definition

An estimator $\hat{\theta}(\mathbf{Z})$ is said to be *minimum variance unbiased estimator* (MVUE) if and only if

C1. $E(\hat{\theta}) = \theta_0$ (unbiased)

C2. $\text{var}(\hat{\theta}) \leq \text{var}(\hat{\theta}_i) \quad \forall i$ satisfying C1 (least variance)

- ▶ The class is restricted to unbiased since biased estimators can always be tuned to have lower variance by sacrificing the bias. Then the comparison becomes difficult.
- ▶ The efficiency of an estimator is used to measure how well it performs relative to an *unbiased* estimator that has the least variance.

Comparing two estimators: Efficiency

Formally, the efficiency of an estimator $\hat{\theta}$ is defined as

$$\text{Efficiency}(\hat{\theta}) = \eta_{\hat{\theta}} = \frac{\text{var}(\hat{\theta}^*)}{\text{var}(\hat{\theta})} \quad (24)$$

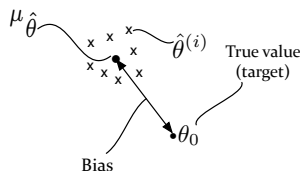
where $\hat{\theta}^*(\mathbf{y})$ has theoretically the lowest variance among all estimators.

Remarks

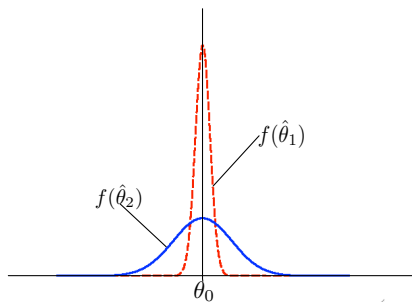
- ▶ The Cramer-Rao's inequality dictates the bound (on achievable variance) and also stipulates the condition under which such an estimator exists.
- ▶ An estimator that achieves this lower bound is said to be the *most efficient* or *fully efficient*.
- ▶ When it is not possible to find an efficient estimator, *relative efficiency* is used.

$$\text{Relative efficiency (\%)} = 100 \times \frac{\sigma_{\hat{\theta}_1}^2}{\sigma_{\hat{\theta}_2}^2} \quad (25)$$

Bias, Variance and Efficiency



Bias is the distance between the center of estimates and the true value, while the variance is a measure of spread around its own center.



The estimator $\hat{\theta}_1$ has lesser spread than $\hat{\theta}_2$, and is therefore relatively more efficient. It produces estimates that have a higher probability (than those of $\hat{\theta}_2$) of being closer to θ_0 .

Most efficient estimator

In seeking the most efficient estimator, it is important to answer the question: **what is the minimum variance achievable by any unbiased estimator?** The celebrated C-R inequality answers this question.

Cramer-Rao inequality

Theorem

Suppose $\hat{\theta}(\mathbf{y})$ is an unbiased estimator of a single parameter θ . Then, if the p.d.f. $f(\mathbf{y}; \theta)$ is regular, the variance of any unbiased estimator is bounded below by $I(\theta)^{-1}$

$$\boxed{\text{var}(\hat{\theta}(\mathbf{y})) \geq (I(\theta))^{-1}} \quad (26)$$

where $I(\theta)$ is the information measure in (6) (or (9)). Further, an estimator $\hat{\theta}^(\mathbf{y})$ that can achieve this lower bound exists if and only if*

$$S(Y_N, \theta) = I(\theta)(\hat{\theta}^*(\mathbf{y}) - \theta) \quad (27)$$

Then, $\hat{\theta}^(\mathbf{y})$ is the **most efficient** estimator of θ .*

Cramer-Rao lower bound

The C-R inequality gives us:

- i. Lowest variance achievable by any unbiased estimator
- ii. Means of deriving that most efficient estimator, if it exists.

The role of Fisher information introduced earlier is clear now.

Larger the information on θ in a given data, lower is the variability and hence the error in $\hat{\theta}$.

Existence of efficient estimator

An alternative form of the condition of existence of an efficient estimator can be given. From (27), the MVUE that achieves the C-R bound exists if and only if

$$\frac{S(Y_N, \theta)}{I(\theta)} + \theta \quad (28)$$

is **independent of θ (sufficiency)** and only dependent on the observations y .

Example: C-R bound

Efficient estimator of mean

Consider the standard problem of estimating the mean of a GWN process $y[k] \sim \mathcal{N}(\mu, \sigma^2)$ from N observations. Find the most efficient estimator of μ .

Solution: Recall from (1)

$$I(\mu) = \frac{N}{\sigma^2} \implies \text{var}(\hat{\mu}) \geq (I(\mu))^{-1} = \frac{\sigma^2}{N} \quad (29)$$

To determine the existence of an estimator that achieves this minimum, construct (28)

$$\frac{S(\mathbf{y}, \theta)}{I(\theta)} + \theta = \frac{\sum_{k=0}^{N-1} (y[k] - \mu)}{N} + \mu = \frac{1}{N} \sum_{k=0}^{N-1} y[k] \quad (30)$$

which is only dependent on the observations \mathbf{y} .

This is none other than the sample mean! In a previous example, we showed that the variance of this estimator is indeed σ^2/N . Thus, we conclude that **the sample mean is the most efficient estimator of the mean of a GWN.**

Existence of an efficient estimator

Whether it is possible to arrive at an efficient estimator depends on two factors:

1. *The parameter θ , or in general, its function $g(\theta)$.* For e.g., in the case of exponentially distributed WN, it turns out that there exists an efficient estimator if $1/\lambda$ is estimated instead of λ .

In parametric modelling, this means that the form of parametrization, i.e., how the parameters enter the model, has an important say in estimation and the estimate.

2. *The probabilistic characteristics of the observed data.* In reality, it is difficult to know the p.d.f. a priori. Then, the existence of an efficient estimator depends on the assumed density function.

Mean Square Error

The mean square error (MSE) of an estimator is its variance with reference to its true value θ_0 .

Definition

The MSE of an estimator is defined as

$$\text{MSE}(\hat{\theta}) = E(\|\hat{\theta} - \theta_0\|_2^2) = E\left(\sum_{i=1}^p (\hat{\theta}_i - \theta_{i0})^2\right) \quad (31)$$

A classical result in estimation relates the bias, variance and MSE.

Theorem

For any estimator $\hat{\theta}$, the following identity holds

$$MSE(\hat{\theta}) = \text{trace}(\Sigma_{\hat{\theta}}) + \|\Delta\hat{\theta}\|_2^2 \quad (32)$$

Proof:

$$\begin{aligned} E(\|\hat{\theta} - \theta_0\|_2^2) &= E(\text{tr}((\hat{\theta} - \theta_0)(\hat{\theta} - \theta_0)^T)) \\ &= \text{tr}(E((\hat{\theta} - \theta_0)(\hat{\theta} - \theta_0)^T)) \\ &= \text{tr}(E((\hat{\theta} - E(\hat{\theta}))(\hat{\theta} - E(\hat{\theta}))^T)) + \text{tr}(E((E(\hat{\theta}) - \theta_0)(E(\hat{\theta}) - \theta_0)^T)) \\ &\quad + 2\text{tr}(E((\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta_0)^T)) \\ &= \text{trace}(\Sigma_{\hat{\theta}}) + \|\Delta\hat{\theta}\|_2^2 \end{aligned}$$

The last identity comes about by recognizing the first term as the trace of $\text{Var}(\hat{\theta})$ and that $E(\hat{\theta} - \theta_0)$ is a deterministic quantity. Consequently the expectation on the second term disappears

$$\text{tr}(E((E(\hat{\theta}) - \theta_0)(E(\hat{\theta}) - \theta_0)^T)) = \text{tr}(\Delta\hat{\theta}\Delta\hat{\theta}^T) = \text{tr}(\Delta\hat{\theta}^T\Delta\hat{\theta}) = \|\Delta\hat{\theta}\|_2^2$$

and the third term vanishes to zero.

- ❶ For unbiased estimators, $\Delta\hat{\theta} = 0$, therefore MSE and $\Sigma_{\hat{\theta}}$ are identical.
- ❷ Since both terms on the RHS of (32) are positive-valued, *estimators that have small MSE naturally require good accuracy and precision.*
- ❸ When $\text{MSE}(\hat{\theta}) \rightarrow 0$ as $N \rightarrow \infty$, the estimator is said to be consistent.
- ❹ It is ideally desirable to build an estimator $\hat{\theta}$ with minimum mean square error. The MMSE problem can be set up by assuming the parameter θ to be a random variable. Therefore, this is useful in a Bayesian estimation framework. The resulting estimator, as it turns out is the conditional expectation $E(\theta|\mathbf{y})$.

Minimum Mean Square Error estimator

Theorem

The MMSE estimator of θ given \mathbf{y} is the conditional expectation

$$\hat{\theta}_{MMSE}(Y) = E(\theta|Y) \quad (33)$$

As in the case of MVUE, the form of MMSE could be non-linear or linear. For practical reasons, linear MMSE estimators are more popular. In fact, when θ and \mathbf{y} follow a joint Gaussian distribution, the linear MMSE is also the optimal MMSE.

Asymptotic bias

Statistical unbiasedness is a desirable property; however, *it is not necessarily the most desirable property*. A *biased* estimator is also considered acceptable provided the bias vanishes for very large samples. For this purpose, asymptotic unbiasedness is defined.

Definition

An estimator is said to be asymptotically unbiased if

$$\lim_{N \rightarrow \infty} \Delta \hat{\theta} = 0 \quad \text{i.e.,} \quad \lim_{N \rightarrow \infty} E(\hat{\theta}) = \theta_0 \quad (34)$$

- ▶ Asymptotic bias is a large sample property. Therefore it is of little interest in situations concerning small samples.

A standard estimator of variance

$$\hat{\sigma}_y^2 = \frac{1}{N} \sum_{k=0}^{N-1} (y[k] - \bar{y})^2 \quad (35)$$

where \bar{y} is the sample mean, *is a biased estimator of σ_y^2 but is asymptotically unbiased.*

- ▶ A *statistically* biased estimator can achieve a variance lower than that of a MVU estimator. However, the variance is no longer a measure for comparing the performance of such estimators since in principle one can shrink the variance to an arbitrarily low (non-zero) value by increasing the bias to a very large value.
- ▶ Thus, a better universal metric is the MSE.

Consistency

An important and desirable large sample property is **consistency**, which examines the convergence of $\hat{\theta}$ to θ_0 as $N \rightarrow \infty$.

An estimator is said to be *consistent* if $\hat{\theta}$ (a RV) converges to θ_0 (a fixed value).

Different forms of consistency arise depending on the notion of convergence one uses:

- ❶ **In probability:** $\hat{\theta}_N \xrightarrow{p} \theta_0$ iff $\lim_{N \rightarrow \infty} \Pr(|\hat{\theta}_N - \theta_0| \geq \varepsilon) = 0, \quad \forall \varepsilon > 0.$
- ❷ **In mean square sense:** $\hat{\theta}_N \xrightarrow{m.s.} \theta_0$ iff $\lim_{N \rightarrow \infty} E((\hat{\theta}_N - \theta_0)^2) = 0.$
- ❸ **Almost sure convergence:** $\hat{\theta}_N \xrightarrow{a.s.} \theta_0$ iff $\hat{\theta}_N \longrightarrow \theta_0$ w.p.1

Order of implication: Almost sure \implies Mean square \implies Probabilistic

Convergence of sequences of random variables

Definition

A sequence of real numbers $\{x_n\}$ is a realization of the sequence of random variables $\{X_n\}$ if x_n is a realization of the RV X_n .

Sequences of RVs on a sample space Ω

$\{X_n\}$ is a sequence of random variables on a sample space Ω if all the RVs belonging to the sequence are mappings from Ω to \mathbb{R} .

One can then have i.i.d or stationary or weakly stationary sequences, etc.

Pointwise convergence

The requirement is that there exist a random variable to which all possible sequences converge on Ω .

Pointwise convergence

Let $\{X_n\}$ be a sequence of random variables defined on a sample space Ω . Then it is pointwise convergent to a random variable X if and only if $\{X_n(\omega)\}$ converges to $X(\omega)$ for all $\omega \in \Omega$. X is called the pointwise limit of the sequence and denoted as

$$\boxed{X_n \xrightarrow{\text{pointwise}} X} \quad (36)$$

Example: PC

Let $\Omega = \{\text{blue}, \text{red}\}$ be the sample space with two sample points. Suppose $\{X_n\}$ is a sequence of RVs such that

$$X_n(\omega) = \begin{cases} \frac{2}{n}, & \omega = \text{blue} \\ 2 + \frac{1}{n}, & \omega = \text{red} \end{cases}$$

Then the sequence converges to a random variable

$$X(\omega) = \begin{cases} 0, & \omega = \text{blue} \\ 2, & \omega = \text{red} \end{cases}$$

Convergence in probability

Idea: The sequence gets very close to a RV X with “high probability”.

Convergence in probability

Let $\{X_n\}$ be a sequence of random variables defined on a sample space Ω and ϵ be a strictly positive number. Then, $\{X_n\}$ is said to be convergent in probability if and only if

$$\lim_{n \rightarrow \infty} \Pr(|X_n - X| > \epsilon) = 0 \quad (37)$$

and denoted by

$$\boxed{X_n \xrightarrow{p} X \quad \text{or} \quad \text{plim}_{n \rightarrow \infty} X_n = X} \quad (38)$$

Example

Consider a sequence of RVs $X_n = \left(1 + \frac{1}{n}\right) X$, where X on $\Omega = \{0, 1\}$ is a discrete RV with p.m.f.

$$p_X(X) = \begin{cases} \frac{1}{5}, & x = 1, \\ \frac{4}{5}, & x = 0 \end{cases}$$

Then $|X_n - X| = 0$ when $X = 0$ (with probability $4/5$) and $|X_n - X| = \frac{1}{n}$ when $X = 1$ (with prob. $1/5$). Therefore,

$$\Pr(|X_n - X| \leq \epsilon) = \begin{cases} \frac{4}{5}, & n < \frac{1}{\epsilon} \\ 1, & n \geq \frac{1}{\epsilon} \end{cases}$$

Mean square convergence

Idea: The sequence gets very close to a RV X in a “squared distance” sense.

Convergence in mean-square

Let $\{X_n\}$ be a sequence of random variables defined on a sample space Ω . Then, $\{X_n\}$ is said to be convergent in mean-square sense if and only if there exists a RV (with finite variance), such that

$$\lim_{n \rightarrow \infty} E[(X_n - X)^2] = 0 \quad (39)$$

i.e., in the sense of a distance metric. The convergence is denoted by

$$\boxed{X_n \xrightarrow{m.s.} X} \quad (40)$$

Example

Consider a sequence of RVs $X_n = \frac{1}{n} \sum_{k=0}^{N-1} x[k]$, where $x[k]$'s are uncorrelated random variables with mean μ and variance σ^2 .

Then, the sequence $\{X_n\}$ converges to a random variable μ in the mean square sense since

$$E((X_n - \mu)^2) = \frac{\sigma^2}{n} \quad \implies \quad \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0 \quad (41)$$

Almost sure convergence

Almost sure convergence is relaxed version of pointwise convergence except that it does accommodate points in Ω where the sequence does not converge.

However, the points $\omega \in \Omega$ on which $\{X_n(\omega)\}$ does not converge pointwise to $X(\omega)$ should be **zero-probability events**.

In other words, **sequences should converge over an interval** (of arbitrarily finite length) in Ω unlike at every point that is required in pointwise convergence.

Then, we write

$$\boxed{X_n \xrightarrow{a.s.} X} \quad (42)$$

Example: a.s. convergence

Consider a sample space $\Omega = [0, 1]$ and a sequence $\{X_n(\omega)\}$ constructed on Ω as

$$X_n(\omega) = \begin{cases} 1, & \omega = 0 \\ \frac{1}{n}, & \omega \neq 0 \end{cases} \quad (43)$$

Examine if the sequence a.s. converges to a (constant) random variable $X(\omega) = 0$.

Deduce that $\lim_{n \rightarrow \infty} X_n(\omega) = \begin{cases} 1, & \omega = 0 \\ 0, & \text{otherwise} \end{cases}$. Note that $\Pr(\omega = 0) = 0$.

Clearly X_n converges to the given random variable $X(\omega)$ except for the event $\omega = 0$, **which is a zero-probability event**. Hence, $X_n(\omega)$ converges to $X(\omega)$ almost surely.

Example 1: Consistency

Sample mean

The sample mean estimator for a WN process has the MSE

$$\text{MSE}(\bar{y}) = \text{var}(\bar{y}) = \frac{\sigma_e^2}{N} \quad (44)$$

This is obviously a m.s. consistent estimator since its $\text{MSE} \rightarrow 0$ as $N \rightarrow \infty$.

Example 1: Consistency

Sample variance

The biased estimator of the variance of a random process was shown to be earlier asymptotically unbiased. For a GWN process with variance σ_e^2 , this estimator is known to have a variance

$$\text{var}(\hat{\sigma}_N^2) = \frac{2(N-1)\sigma_e^4}{N^2} \quad (45)$$

Therefore, it is mean-square consistent.

Remarks

- ▶ Consistency essentially guarantees that increasing the number of observations takes us “closer” to the true value. Therefore, it is practically one of the most important properties of an estimator.
- ▶ There are several estimators that are not consistent. A popular one is the *periodogram*, which estimates the power spectral density of a signal.
- ▶ For biased estimators, mean square consistency also implies asymptotic unbiasedness because

$$\text{MSE}(\hat{\theta}) = \text{bias}^2 + \text{var}(\hat{\theta})$$

Running summary

To recap the key points until now:

- ▶ The goodness (accuracy, precision, etc.) of an estimate depends on two factors: (i) information content in the data and (ii) properties of the estimator
 - ▶ Information content is measured by Fisher's information, which is based on the likelihood function. It is a measure of the quality of a given dataset w.r.t. estimating θ and is regardless of the form of $\hat{\theta}$.
- ▶ Six properties of an estimator are usually important: *bias, variance, efficiency, mean square error, asymptotic bias and consistency*.
 - ▶ Efficiency and consistency are the two most important criteria

To recap the key points until now:

- ▶ C-R inequality gives us the lowest variability (error) that can be achieved by an unbiased estimator. The bound is the inverse of FI.
- ▶ Whenever it becomes difficult to estimate or design a 100% efficient estimator or even a MVUE, a **best linear unbiased estimator** is sought.
- ▶ Consistency guarantees convergence of the estimate to the true value

Motivation

Up to this point we studied metrics for quantifying the goodness of data and estimators. Now we raise an important question:

Given an observation vector \mathbf{y} and a *point* estimate $\hat{\theta}$ what can be said about the true value θ_0 ?

Interval estimates and hypothesis testing

Two related problems are:

- ➊ **Confidence intervals:** What is the interval in which the true value resides? Only intervals are sought since the true value cannot be estimated precisely.
- ➋ **Hypothesis testing:** Given $\hat{\theta}$ how do we test claims on the true parameters?

To be able to answer the above questions it is necessary to determine the probability distribution of an estimate.

Introductory remarks

The distribution of estimate generally depends on three factors:

- 1 *Randomness in observations*: It is a crucial factor since it is the “feed” to the estimator. It is the source of uncertainty in estimate.
- 2 *Form of estimator*: When the estimator is linear (e.g., sample mean, BLUE estimator) the transformation of the $f(\mathbf{y}; \boldsymbol{\theta})$ can be easily studied. Non-linear estimators naturally pose a challenge, except under very special conditions.
- 3 *Sample size*: A large body of estimation literature is built on the large sample size assumption. Small sample sizes not only affect the distribution but also the consistency property of an estimator!

Multiplication by \sqrt{N}

Distributions are quite often stated for $\sqrt{N}(\hat{\theta} - \theta_0)$ (or at times $\sqrt{N}\hat{\theta}_N$) instead of $\hat{\theta}_N$ itself. This is because asymptotically $(\hat{\theta}_N - \theta_0)$ converges to a constant (mostly zero), whereas $\sqrt{N}(\hat{\theta} - \theta_0)$ converges to a random variable with a meaningful distribution.

Example

From the previous sections, we know that the sample mean is a consistent estimator of the mean. This means $\bar{y}_N - \mu_y$ converges to zero as $N \rightarrow \infty$. On the other hand, $\sqrt{N}(\bar{y}_N - \mu_y)$ converges to a random variable with mean zero and finite variance.

$$E(\hat{\theta}_N) = \mu_y; \quad \text{var}(\bar{y}) = \frac{\sigma_y^2}{N} \implies (\bar{y} - \mu_y) \xrightarrow{m.s.} 0$$

but, $\sqrt{N}(\bar{y} - \mu_y) \xrightarrow{m.s.} \sigma_y^2$

Convergence in distribution

In order to study the asymptotic distributional properties of an estimator, it is necessary to first understand the notion of **convergence in distribution** of a sequence of RVs.

Definition

A sequence of random variables $\{X_N\}$, each possessing a distribution function $F(x_N)$ **converges in distribution** if the sequence of those distributions $\{F(x_N)\}$ (sometimes written as $F_N(x)$) converges to a distribution function $F(x)$. The random variable X associated with $F(x)$ is said to be the *limit in distribution* of the sequence, indicated as

$$X_n \xrightarrow{d} X \quad (46)$$

Note that the theorem speaks of convergence of distributions, not the RVs themselves.

Central Limit Theorem

The CLT is one of the most celebrated and landmark results in estimation theory. Historically it is nearly seven decades old and has undergone several modifications. The basic version due to Lindeberg and Levy is as follows.

Theorem (Central Limit Theorem)

The uniformly weighted sum of N i.i.d. random variables $\{X_n, n = 1, \dots, N\}$

$$\bar{X} = \sum_{n=1}^N \frac{X_n}{N}, \quad \text{where} \quad E[X_n] = \mu < \infty, \quad \text{var}(X_n) = \sigma^2 < \infty \quad (47)$$

converges in distribution as

$$\sqrt{N} \left(\frac{\bar{X} - \mu}{\sigma} \right) \xrightarrow{d} \mathcal{N}(0, 1) \quad (48)$$

- ▶ The conditions of independence and identical distributions are not heavily restrictive. Versions of CLT which place some minor additional requirements on the moments or the autocorrelation functions are also available.
- ▶ Generalizations and extensions to other random objects such as matrices and polytopes are available.
- ▶ Note that the sum is none other than the sample mean of the N random variables.
- ▶ It is also conventional to state that the distribution of \bar{X} is *asymptotically normal* or simply write as

$$\sqrt{N}(\bar{X} - \mu) \sim \mathcal{AN}(0, \sigma^2) \quad (49)$$

Limitations

- ▶ The CLT provides us with a tool for deriving distributions of parameter estimates from linear estimators with known distributional properties of the data. Many standard results on distributions of estimators such as sample mean, sample variance, linear least squares estimates can be derived through CLT.
- ▶ However, when the estimator is a complicated function of the observations, further simplifying approximations or the use of modern (Monte-Carlo) methods such as bootstrapping or surrogate data analysis have to be employed.

Confidence regions

The term “confidence region” essentially refers to the interval containing the true value with $< 100\%$ confidence. Ideally one would like to have a narrow interval with maximum confidence. However, these are conflicting requirements because a higher degree of confidence is associated with a wider band.

The procedure for constructing a confidence interval is a two-step process.

Step 1: Construct a probabilistic interval for the error $\hat{\theta}_N - \theta_0$ using the knowledge of the distribution (or density) of $\hat{\theta}_N$, the bias and variance properties and the specified degree of confidence $100(1 - \alpha)\%$.

Step 2: Convert this *probabilistic* interval into a *confidence* region for θ_0 by an algebraic manipulation.

Confidence interval for mean

Assume that the sample mean \bar{y} is used as an estimator of the mean μ_y from a single record of data.

Goal: To obtain a confidence region for μ_y

Assume that σ_y^2 is known. Invoking CLT, $\sqrt{N} \left(\frac{\bar{y} - \mu_y}{\sigma_y} \right) \sim \mathcal{N}(0, 1)$

From the properties of a Gaussian distribution,

$$-1.96 \leq \sqrt{N} \frac{\bar{y} - \mu_y}{\sigma_y} \leq 1.96 \quad (\text{with 95\% probability})$$

$$\Rightarrow \mu_y \in \left[\bar{y} - \frac{1.96}{\sqrt{N}} \sigma_y, \bar{y} + \frac{1.96}{\sqrt{N}} \sigma_y \right] \quad (\text{with 95\% confidence}) \quad (50)$$

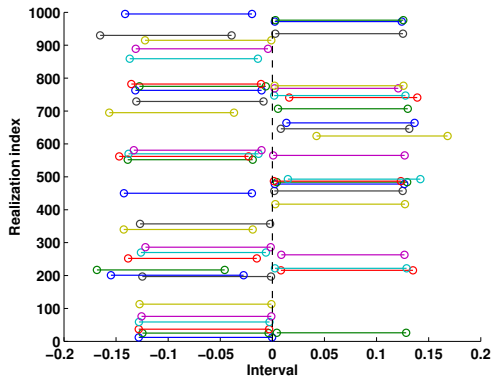
The $100(1 - \alpha)\%$ CI for the mean is obtained by replacing 1.96 with ζ_c such that $\Pr(\zeta > \zeta_c) = \alpha/2$ (using the standard Gaussian distribution).

Interpretation

The confidence interval (CI) should be interpreted with care. Consider the case of a 95% CI for mean. Suppose that we have 1000 records of data, from each of which we can obtain an estimate $\bar{y}^{(i)}$, $i = 1, \dots, 1000$, from each of which a 95% C.I. can be constructed. Then, out of 1000 such CIs, roughly 950 intervals would have correctly captured the true mean.

Simulation example

In one such simulation study, it turns out that 51 intervals do not contain the true value $\mu_0 = 0$ as shown below.



Remarks

- ▶ The width of the CI is only dependent on the standard error in the estimate σ_y/\sqrt{N} . In general, **the width depends on the variability of the process and the sample size** (for a consistent estimator)
- ▶ Narrower the width of the interval at a fixed α , better is the estimator. A consistent estimator produces zero-width CI asymptotically.
- ▶ For correlated processes, the CI has to be re-derived because $\text{var}(\bar{y})$ is influenced by the correlation structure.

Confidence intervals

- ① **Mean:** Small sample, variance unknown.

$$\mu_y \in [\bar{y} - t_{\alpha/2}(N-1)\hat{\sigma}_y, \bar{y} + t_{\alpha/2}(N-1)\hat{\sigma}_y] \quad (\text{with 95\% confidence}) \quad (51)$$

- ② **Variance:** Gaussian population, random samples

$$\frac{(n-1)S^2}{\chi_{\alpha/2, n-1}^2} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi_{1-\alpha/2, n-1}^2} \quad (52)$$

One-sided: Lower and upper confidence bounds

$$\frac{(n-1)S^2}{\chi_{\alpha, n-1}^2} \leq \sigma^2, \quad \sigma^2 \leq \frac{(n-1)S^2}{\chi_{1-\alpha, n-1}^2} \quad (53)$$

Several standard texts on statistics present the theory of hypothesis testing and confidence interval construction (see Johnson, 2011; Ogunnaike, 2010).

For non-linear estimators, modern empirical methods are used to obtain the distributions of estimates via the generation of **surrogate data** or *pseudo-population*. Monte-Carlo simulations and bootstrapping methods are increasingly being used for this purpose.

Hypothesis testing

Once the distributional properties of an estimator are known, it is possible to answer the two questions raised earlier, i.e., pertaining to hypothesis testing and confidence interval construction. Both are in fact related problems.

Hypothesis testing involves a statistical test for a claim made by the analyst with regards to the properties of the process of interest or model parameters using the observations as an evidence.

Examples:

- ▶ Average temperature of a reactor is at a specified value.
- ▶ Model parameters are truly zero.
- ▶ The given series is white (unpredictable)

Procedure for hypothesis testing

A hypothesis test typically consists of the following steps.

- 1 *Formulate the null hypothesis H_0* based on the postulate or the claim. Choose an appropriate alternate hypothesis H_a .
- 2 *Choose an appropriate statistic ζ* for the test. The statistic is generally a linear or non-linear function of the parameter(s) involved in the hypothesis.
- 3 *Compute the test statistic* from the given observations. Denote this by ζ_o .
- 4 *Make a decision.* Retain or discard the null hypothesis by applying a certain criterion to the observed statistic (three different approaches).

No hypothesis test can result in a perfect decision!

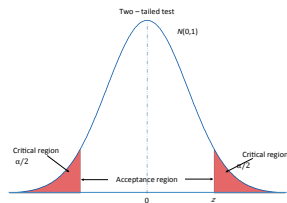
Errors in hypothesis testing

Any hypothesis test is marred by two errors - Type I and II errors. Typically, the first type, known as the α risk or the **significance level** is specified.

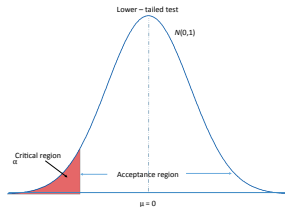
Decision \longrightarrow Truth \downarrow	Fail to Reject H_0	Reject H_0
H_0 True	Correct Decision Probability: $1 - \alpha$	Type I Error Probability (Risk): α
H_a True	Type II Error Risk: β	Correct Decision Probability: $1 - \beta$

One of the two errors has to be specified for making a decision in hypothesis testing. **It is a common practice to specify the α risk.**

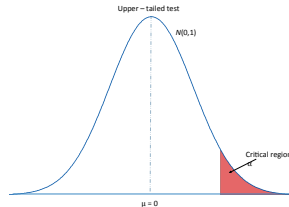
Graphical understanding: One sample test for mean



(a) Two tailed test



(b) Lower-tailed test



(c) Upper-tailed test

The α risk (probability of making Type I error) depends on the type of **alternate hypothesis** and the **sampling distribution** of the statistic.

Decision making in hypothesis testing

There are **three** different approaches to making decisions in hypothesis testing, *all of which lead to the same result*.

- ➊ **Critical value approach:** Determine a critical value (for a given risk) and compare the observed statistic against it.
- ➋ **p -value method:** Determine the probability of obtaining a value more extreme than the observed and compare this probability against a user-specified value (risk).
- ➌ **Confidence interval approach:** Construct the confidence region (for a given risk) and determine if the postulated value falls within the region.

p-value

The *p*-value is the probability of observing a more extreme value of the statistic than the observed value.

- ▶ It is computed under the null hypothesis being held to be true.
- ▶ The sign or the direction of the extreme value depends on the alternate hypothesis just as the way the critical value does.
- ▶ The *p*-value computation is fairly straightforward. For example, in an upper tail test, the *p*-value is the $\Pr(\zeta \geq \zeta_o)$. If the *p*-value $\leq \alpha$, then H_0 is rejected.

Example: Hypothesis testing

An engineer measures the (controlled) temperature of a reactor over a period of 3 hours at a sampling interval of $T_s = 15$ sec. The sample average of the $N = 720$ readings is calculated to be $\bar{y} = 90.1826^\circ\text{C}$. Based on this observation, the engineer claims that the temperature is at its set-point $T_0 = 90^\circ\text{C}$ on the average.

To test this claim, the formal hypothesis test is set up as follows.

$$H_0 : \mu_y = 90 \quad H_a : \mu_y \neq 90 \quad \text{two-tailed test}$$

Example: Hypothesis testing

... contd.

Assume that the temperature series has white-noise characteristics. Then we know that for the large sample case,

$$\sqrt{N} \left(\frac{\bar{Y} - \mu}{\sigma} \right) \xrightarrow{d} \mathcal{N}(0, 1)$$

An appropriate test statistic suited for the purpose is therefore,

$$Z = \frac{\bar{y} - \mu_0}{\sigma / \sqrt{N}} \quad (54)$$

where μ_0 is the true value assumed in H_0 . For the example, $\mu_0 = T_0 = 90^\circ\text{C}$. Assume that σ is known to be 2°C . Then the observed statistic is $z_0 = 2.45$.

Decision making

- i. **Critical value approach:** The critical value at $\alpha = 0.05$ is $z_c = 1.96$. Since $z_o > z_c$, the null hypothesis is rejected.
- ii. **p-value approach:** The $\Pr(|Z| > z_0 = 2.45) = 0.0143 < \alpha = 0.05$. Hence H_0 stands rejected in favour of H_a .
- iii. **C.I. approach:** The $100(1 - \alpha)$ C.I. for the average temp. is $(90.0365, 90.3287)$, which does not include the postulated value. Hence H_0 stands rejected in favour of H_a .

On the average the temperature is not at its set-point, i.e., the engineer's claim that $H_0 : \mu = 90^\circ$ (set-point), stands rejected in favour of the alternate hypothesis.

Remarks

- ▶ How would you adjust the significance level just enough so that H_0 is not rejected?
- ▶ **Note:** The assumption of known σ can be relaxed. It can be estimated σ from data, in which case ζ has a t -distribution with $\nu = N - 1$ degrees of freedom.

Confidence intervals in Hypothesis testing

In a general situation, the C.I. approach for testing hypothesis (all three forms) is as follows :

Procedure

- ❶ Specify the significance level.
- ❷ Depending on the hypothesis, construct the appropriate C.I. and apply the test
 - ❶ **Two-sided:** $100(1 - \alpha/2)\%$ C.I. If postulated value is not within the C.I., reject H_0 .
 - ❷ **One-sided:** Reject H_0 if the postulated value is greater or lesser than the bound, for the upper- and lower-tailed test, respectively.

Summary

- ▶ The end goal of an estimation exercise is to arrive at an interval estimate or a confidence region for the true value
- ▶ Distributions of estimates provide the necessary information to move from a point to an interval estimate
- ▶ The distributional properties also facilitate hypothesis testing, which involves a systematic and statistical way of testing the claims related to a process and/or a model
- ▶ Arriving at distributions is non-trivial, but the CLT comes to the rescue for linear estimators
- ▶ The case of non-linear estimation is more complicated and calls for the use of modern tools such as Monte-Carlo simulations and Bootstrapping methods.

Bibliography I

- Fisher, R. (1922). On the mathematical foundations of theoretical statistics. *Philosophical Transactions Royal Society London A*, 222, pp. 309–368.
- (1950). *Contributions to mathematical statistics*. New York, USA: John Wiley & Sons, Inc.
- Garthwaite, P., I. Jolliffe, and B. Jones (2002). *Statistical Inference*. New York, USA: Oxford University Press.
- Johnson, R. A. (2011). *Miller and Freund's: Probability and Statistics for Engineers*. Upper Saddle River, NJ, USA: Prentice Hall.
- Ogunnaike, B. A. (2010). *Random Phenomena: Fundamentals of Probability and Statistics for Engineers*. Boca Raton, FL, USA: CRC Press, Taylor & Francis Group.
- Tangirala, A. K. (2014). *Principles of System Identification: Theory and Practice*. Boca Raton, FL, USA: CRC Press, Taylor & Francis Group.